

# PARALLEL ADAPTIVE IMPORTANCE SAMPLING

COLIN COTTER\*, SIMON COTTER†, AND PAUL RUSSELL‡

**Abstract.** Markov chain Monte Carlo methods are a powerful and commonly used family of numerical methods for sampling from complex probability distributions. As applications of these methods increase in size and complexity, the need for efficient methods which can exploit the parallel architectures which are prevalent in high performance computing increases. In this paper, we aim to develop a framework for scalable parallel MCMC algorithms. At each iteration, an importance sampling proposal distribution is formed using the current states of all of the chains within an ensemble. Once weighted samples have been produced from this, a state-of-the-art resampling method is then used to create an evenly weighted sample ready for the next iteration. We demonstrate that this parallel adaptive importance sampling (PAIS) method outperforms naive parallelisation of serial MCMC methods using the same number of ensemble members, for low dimensional problems, and in fact shows better than linear improvements in convergence rates with respect to the number of ensemble members. We also introduce a new resampling strategy, approximate multinomial resampling (AMR), which while not as accurate as other schemes is substantially less costly for large ensemble sizes, which can then be used in conjunction with PAIS for complex problems.

**Key words.** MCMC, parallel, importance sampling, Bayesian, inverse problems.

**1. Introduction.** Markov chain Monte Carlo (MCMC) methods are a powerful family of tools that allow us to sample from complex probability distributions. MCMC methods were first developed in the 70s [12], and with the development of faster, more powerful computers, have become ever more important in a whole range of fields in statistics, science and engineering. In particular, when considering Bayesian inverse problems, each MCMC step may involve the numerical solution of one or more PDE. As many samples are usually required before Monte Carlo error is reduced to acceptable levels, the application of MCMC methods to these types of problem remain frustratingly out of our grasp.

Many advances have been made in the field of MCMC to design ever more complex methods that propose moves more intelligently, leading to rapidly converging approximations. Function space versions of standard methods such as the random walk Metropolis-Hastings (RWMH) algorithm or the Metropolis adjusted Langevin algorithm (MALA), whose convergence rates are independent of dimension have been developed [6]. The hybrid (or Hamiltonian) Monte Carlo (HMC) method uses Hamiltonian dynamics in order to propose and accept moves to states which are a long way away from the current position [26], and function space analogues of this have also been proposed [1]. Riemann manifold Monte Carlo methods exploit the Riemann geometry of the parameter space, and are able to take advantage of the local structure of the target density to produce more efficient MCMC proposals [10]. This methodology has been successfully applied to MALA-type proposals and methods which exploit even higher order gradient information [2]. These methods allow us to explore the posterior distribution more fully with fewer iterations.

Simultaneously, great strides are continually being made in the development of computing hardware. Moore's law, which predicted that the number of transistors that can fit onto a single microchip will double every two years, has been largely followed since the early 70s [17]. In recent times, it has become necessary to use parallel

\*Department of Mathematics, Imperial College, London, UK

†School of Mathematics, University of Manchester, Manchester, UK. e: simon.cotter@manchester.ac.uk. SLC is grateful for EPSRC First grant award EP/L023393/1

‡School of Mathematics, University of Manchester, Manchester, UK

architectures in order for this trend to continue. The efficient exploitation of these architectures is the key to solving many of the computational challenges that we currently face.

As such, the development of efficient parallel MCMC algorithms is an important area for research. Since MCMC methods can be naively parallelised by simply running many independent chains in parallel, the focus needs to be on the development of methods which gain some added benefit through parallelism. One class of parallel MCMC method uses multiple proposals, with only one of these proposals being accepted. Examples of this approach include multiple try MCMC [16] and ensemble MCMC [18]. In [3], a general construction for the parallelisation of MCMC methods was presented, which demonstrated speed ups of up to two orders of magnitude when compared with serial methods.

In this paper, we present a framework for parallelisation of importance sampling, which can be built around many of the current Metropolis-based methodologies in order to create an efficient target proposal from the current state of all of the chains in the ensemble. The idea is to consider the current state on each of a set of parallel chains as an ensemble, and to resample using a transformation based on optimal transport. Samplers based on optimal transport have also been considered in [7].

In Section 2 we introduce some mathematical preliminaries upon which we will later rely. In Section 3 we present the general framework of the PAIS algorithm. In Section 4 we consider adaptive versions of PAIS which automatically tune algorithmic parameters concerned with the proposal distributions. In Section 5 we introduce the approximate multinomial resampling (AMR) algorithm which is a less accurate but faster alternative to resamplers which solve the optimal transport problem exactly. In Section 6 we present some numerical examples, before a brief conclusion and discussion in Section 7.

**2. Preliminaries.** In this Section we will introduce preliminary topics and algorithms that will be referred to throughout the paper.

**2.1. Bayesian inverse problems.** In this paper, we focus on the use of MCMC methods for characterising posterior probability distributions arising from Bayesian inverse problems. We wish to learn about a particular unknown quantity  $u$ , of which we are able to make direct or indirect noisy observations. For now we say that  $x$  is a member of a Hilbert space  $X$ .

The parameter  $x$  is observed through the observation operator  $\mathcal{G} : X \rightarrow \mathbb{R}^d$ . Since observations are never perfect, we assume that these measurements  $D$  are subject to Gaussian noise, so that

$$D = \mathcal{G}(x) + \varepsilon, \quad \varepsilon \sim \mu_\varepsilon = \mathcal{N}(0, \Sigma). \quad (2.1)$$

For example, if  $x$  are the rates of reactions in a chemical system,  $\mathcal{G}$  might return the quantities of each chemical species at a particular time, or some summary of this information.

These modelling assumptions allow us to construct the likelihood of observing the data  $D$  given the parameter  $x = x^*$ . Rearranging (2.1) and using the distribution of  $\varepsilon$ , we get:

$$\mathbb{P}(D|x = x^*) \propto \exp\left(-\frac{1}{2}\|\mathcal{G}(x^*) - D\|_\Sigma^2\right) = \exp(-\Phi(x^*)), \quad (2.2)$$

where  $\|y_1 - y_2\|_\Sigma$  is the Mahalanobis distance between  $y_1$  and  $y_2 \in \mathbb{R}^d$ .

As discussed in [5,27], in order for this inverse problem to be well-posed in the Bayesian sense, we require the posterior distribution,  $\mu_Y$ , to be absolutely continuous with respect to the prior,  $\mu_0$ . A minimal regularity prior can be chosen informed by regularity results of the observational operator  $\mathcal{G}$ . Given such a prior, then the Radon-Nikodym derivative of the posterior measure,  $\mu_Y$ , with respect to the prior measure,  $\mu_0$ , is proportional to the likelihood:

$$\frac{d\mu_Y}{d\mu_0} \propto \exp(-\Phi(x^*)). \quad (2.3)$$

**2.2. Particle filters and resamplers.** In several applications, data must be assimilated in an “online” fashion, with up to date observations of the studied system being made available on a regular basis. In these contexts, such as in weather forecasting or oceanography, data is incorporated using a filtering methodology. One popular filtering method is the particle filter, the first of which was dubbed the Bootstrap filter [11]. In this method, a set of weighted particles is used to represent the posterior distribution. The positions of the particles are updated using the model dynamics. Then, when more observations are made available, the relative weights of the particles are updated to take account of this data, using Bayes’ formula. Other filtering methods, such as the Kalman filter [14] and ensemble Kalman filter [8], have also been developed which are often used within the data assimilation community.

One advantage of the particle filter is that there are convergence results for this method as the number of particles is increased. One downside is that the required ensemble size increases quickly with dimension, making it difficult to use in high-dimensional problems. Another downside is that the effective sample size decreases at each iteration, resulting in degeneration of the approximation of the posterior. One way to tackle this is to employ a resampling scheme. The aim of a successful resample is to take the unevenly weighted ensemble and return a new ensemble of particles with even weights which is highly correlated to the original samples.

The ensemble transform particle filter (ETPF) proposed by Reich [20] makes use of optimal transportation as described in [28,29]. The transform takes a sample of weighted particles  $\{y_i\}_{i=1}^M$  from  $\mu_Y$  and converts it into a sample of evenly weighted particles  $\{x_i\}_{i=1}^M$  from  $\mu_X$ , by means of defining a coupling  $T^*$  between  $Y$  and  $X$ . Given that a trivial coupling  $T^t$  always exists in the space of transference plans,  $\Pi(\mu_X, \mu_Y)$ , we can find a coupling  $T^*$  which maximises the correlation between  $X$  and  $Y$  [4]. This coupling is the solution to a linear programming problem in  $M^2$  variables with  $2M - 1$  constraints. Maximising the correlation ensures that the new sample is as much like the original sample as possible with the additional property that the sample is evenly weighted.

A Monte Carlo algorithm can be implemented to resample from a weighted ensemble. We create a weighted sample, then solve the optimal transport problem which produces the coupling described above, we can draw a new sample from the evenly weighted distribution. Reich suggests using the mean of the evenly weighted distribution to produce a consistent estimator. Analysis of this method shows that as the ensemble size increases, the statistics of the evenly weighted sample approach those of the posterior distribution.

**2.3. Deficiencies of Metropolis-type MCMC schemes.** All MCMC methods are naively parallelisable. One can take a method and simply implement it simultaneously over a set of processors in an ensemble. All of the states of all of the

ensemble member can be recorded, and in the time that it takes one MCMC chain to draw  $N$  samples,  $M$  ensemble members can draw  $NM$  samples.

However, we argue that this is not an optimal scenario. First of all, unless we have a lot of information about the posterior, we will initialise the algorithm’s initial state in the tails of the distribution. The samples that are initially made as the algorithm finds its way to the mode(s) of the distribution cannot be considered to be samples from the target distribution, and must be thrown away. This process is known as the burn-in. In a naively parallelised scenario, each ensemble member must perform this process independently, and therefore mass parallelisation makes no inroads to cutting this cost.

Moreover, many MCMC algorithms suffer from poor mixing, especially in multimodal systems. The amount of samples that it takes for an MCMC trajectory to switch between modes can be large, and given that a large number of switches are required before we have a good idea of the relative probability densities of these different regions, it can be prohibitively expensive.

Another aspect of Metropolis-type samplers is that information computed about a proposed state is simply lost if we choose to reject that proposal in the Metropolis step. An advantage of importance samplers is that no evaluations of  $\mathcal{G}$  are ever wasted since all samples are saved along with their relative weighting.

Moreover, a naively parallelised MCMC scheme is exactly that - naive. Intuition suggests that we can gain some speed up by sharing information across the ensemble members, and this is what we wish to demonstrate in this paper.

These deficiencies of the naive method of parallelising MCMC methods motivated the development of the Parallel Adaptive Importance Sampler (PAIS). In the next section we will introduce the method in its most general form.

**3. The Parallel Adaptive Importance Sampler (PAIS).** Importance sampling can be a very efficient method for sampling from a probability distribution. A proposal density is chosen, from which we can draw samples. Each sample is assigned a weight given by the ratio of the target density and the proposal density at that point. They are efficient when the proposal density is concentrated in similar areas to the target density, and incredibly inefficient when this is not the case. The aim of the PAIS is to use an ensemble of states, each coming from a MCMC chain, to construct a proposal distribution which will be as close as possible to the target density. If this ensemble is large enough, the distribution of states will be representative of the target density.

The proposal distribution could be constructed in many different ways, but we choose to use a mixture distribution, made up of a sum of MCMC proposal distributions for each of the members of the ensemble. Once the proposal is constructed, we can sample a new set of states from the proposal distribution, and each is assigned a weight given by the ratio of the target density and the proposal mixture distribution density. Assuming that our proposal distribution is a good one, then the variance of the weights will be small, and we will have many useful samples. Finally, we need to create a set of evenly weighted samples which best represent this set of weighted samples. This is achieved by implementing a resampling algorithm. Initially we will use the ETPF algorithm [20], although we will suggest an alternative strategy in Section 5. The output of the resampling algorithm gives us a set of evenly weighted samples that we believe represents the target distribution well, and from this point we can iterate the process once again. The algorithm is summarised in Table 3.1.

We wish to sample states  $x \in X$  from a posterior probability distribution  $\mu_D$ , where

```

 $\mathbf{X}^{(0)} = \mathbf{X}_0 = [x_1^{(0)}, x_2^{(0)}, \dots, x_M^{(0)}]^T$ 
for  $i = 0, 1, 2, \dots, N$  do
   $\mathbf{Y}^{(i)} = [y_1^{(i)}, y_2^{(i)}, \dots, y_M^{(i)}]^T, \quad y_j^{(i)} \sim \nu(\cdot; x_j^{(i)})$ 
   $\chi(y; \mathbf{X}^{(i)}) = \frac{1}{M} \sum_{j=1}^M \nu(y; x_j^{(i)})$ .
   $\mathbf{W}^{(i)} = [w_1^{(i)}, w_2^{(i)}, \dots, w_M^{(i)}]^T, \quad w_j^{(i)} = \frac{\pi(y_j^{(i)})}{\chi(y_j^{(i)}; \mathbf{X}^{(i)})}$ .
  Resample:  $(\mathbf{W}^{(i)}, \mathbf{Y}^{(i)}) \rightarrow (\frac{1}{M} \mathbf{1}, \mathbf{X}^{(i+1)})$ 
end for

```

TABLE 3.1

A pseudo-code representation of the Parallel Adaptive Importance Sampler (PAIS).

$D$  represents our data which we wish to use to infer  $x$ . Since we have  $M$  ensemble members, we represent the current state of all of the Markov chains as a vector  $\mathbf{X} = [x_1, x_2, \dots, x_M]^T$ . We are also given a transition kernel  $\nu(\cdot, \cdot)$ , which might come from an MCMC method, for example the random walk Metropolis-Hastings proposal density  $\nu(\cdot, x) \sim \mathcal{N}(x, \beta^2)$ , where  $\beta^2 \in \mathbb{R}$  defines the variance of the proposal.

Since the resampling does not give us a statistically identical sample to that which is inputted, we cannot assume that the samples  $\mathbf{X}^{(i)}$  are samples from the posterior. Therefore, as with serial importance samplers, the weighted samples  $(\mathbf{W}^{(i)}, \mathbf{Y}^{(i)})_{i=1}^N$  are the samples from the posterior that we will analyse.

The key is to choose a suitable transition kernel  $\nu$  such that if  $X^{(i)}$  is a good representative sample of the posterior, then the mixture density  $\chi(\cdot; \mathbf{X}^{(i)})$  is a good approximation of the posterior distribution. If this is the case, the newly proposed states  $\mathbf{Y}^{(i)}$  will also be a good sample of the posterior with low variance in the weights  $\mathbf{W}^{(i)}$ .

In Section 6, we will demonstrate how the algorithm performs, primarily using RWMH proposals. We do not claim that this choice is optimal, but is simply chosen as an example to show that sharing information across ensemble members can improve on the original MCMC algorithm and lead to convergence in fewer evaluations of  $\mathcal{G}$ . This is important since if the inverse problem being tackled involves computing the likelihood from a very large data set this could lead to a large saving of computational cost. We have observed that using more complex (and hence more expensive) kernels  $\nu$ , does not significantly improve the speed of convergence of the algorithm for simple examples [25].

Care needs to be taken when choosing the proposal distribution to ensure that the proposals are absolutely continuous with respect to the posterior distribution. In Section 6.3 we will consider an inverse problem with Gamma priors. Since the posterior distribution in this context is heavy-tailed, if we use a lighter tailed distribution for  $\nu$ , such as a Gaussian kernel, then importance weights in the tails will be unbounded, which will hamper convergence of the algorithm.

**4. Automated tuning of Algorithm Parameters.** Efficient selection of scaling parameters in MCMC algorithms is critical to achieving optimal mixing rates and hence achieving fast convergence to the target density. One aspect worthy of consideration with the PAIS, is finding an appropriate proposal kernel  $\nu$  such that the mixture distribution  $\chi$  is a close approximation to the posterior density  $\pi$ . If the proposal distribution is too over-dispersed, then the algorithm will often propose states in

the tails of the distribution, resulting in larger variance of the weights, and therefore slower convergence to the posterior distribution. Similarly, if the proposal distribution is under-dispersed, the proposals will be highly correlated with the previous states, and the algorithm will take a long time to fully explore the parameter space, and worse, will lead regularly to states proposed in the tails of the proposal distribution with very large weights. It is therefore necessary to find a proposal distribution which is slightly over-dispersed to ensure the entire posterior is explored [9], but is as close to the posterior as possible.

Most commonly used MCMC proposals have parametric dependence which allows the user to control their variance. For example, in the RWMH proposal  $y = x + \beta\eta$ , the parameter  $\beta$  controls how correlated the proposal state  $y$  is to the current state  $x$ . Therefore the proposal distributions can be tuned such that they are slightly over-dispersed, as described above. This tuning can take place during the burn-in phase of the algorithm. Algorithms which use this method to find optimal proposal distributions are known as adaptive MCMC algorithms, and have been shown to be convergent provided that they satisfy certain conditions [22, 23].

Algorithms which use mixture proposals, e.g. PAIS, must tune the variance of the individual kernels within the proposal mixture. This adaptivity during the burn-in has some added benefits over and above finding an optimal parameter regime for the algorithm. If the initial value of the proposal variances is chosen to be very large, then proposed moves will be made far and wide, expediting the early mode-finding stages of the algorithm. Adaptively reducing the proposal variances to an optimal value then allows us to explore each region efficiently. The fact that we have an ensemble of chains allows us to assess quickly and effectively what the optimal variance of the proposal distributions should be.

The alternative to using adaptive procedures to tune the scaling parameters is to perform exploratory simulations to find the optimal regimes by trial and error. This can be very costly and the optimal parameters cannot realistically be found to more than a couple of significant figures. It is therefore important that MCMC algorithms provide a feasible adaptive strategy.

In many MCMC algorithms such as the Random Walk Metropolis-Hastings (RWMH) algorithm, the optimal scaling parameter can be found by searching for the parameter value which gives an optimal acceptance rate, e.g. for near Gaussian targets the optimal rates are 23.4% for RWMH and 57.4% for MALA [21]. Unlike Metropolis-Hastings algorithms, the PAIS algorithm does not accept or reject proposed values, so we need another method of measuring the optimality of  $\beta$ . Section 4.1 gives some possible methods for tuning  $\beta$ .

#### 4.1. Statistics for Determining the Optimal Scaling Parameter.

**4.1.1. Determining optimal scaling parameter using error analysis.** MCMC algorithms can be assessed by comparing their approximation of the posterior to the analytic distribution, in cases where the posterior distribution can be computed using alternative methods. To assess this, a distance metric on distributions must be chosen. Examples are the relative error between the sample moments and the posterior's moments, or the relative  $L^2$  error between the true density,  $\pi(x|D)$ , and the constructed histogram. The relative error in the  $m$ -th moment is given by:

$$\left| \frac{N^{-1} \sum_{i=1}^N x_i^m - \mathbb{E}[X^m]}{\mathbb{E}[X^m]} \right|, \quad (4.1)$$

where  $\{x_i\}_{i=1}^N$  is a sample of size  $N$ . The relative  $L^2$  error between a continuous function to a piecewise constant function,  $e$ , can be given by considering the difference in mass between the normalised histogram of the samples and the posterior distribution over a set of disjoint sets or “bins”:

$$e^2 = \sum_{i=1}^{n_b} \left[ \int_{R_i} \pi(a|D) da - vB_i \right]^2 / \sum_{i=1}^{n_b} \left[ \int_{R_i} \pi(a|D) da \right]^2, \quad (4.2)$$

where the regions  $\{R_i\}_{i=1}^{n_b}$  are the  $d$ -dimensional histogram bins, so that  $\bigcup_i R_i \subseteq X$  and  $R_i \cap R_j = \emptyset$ ,  $n_b$  is the number of bins,  $v$  is the volume of each bin, and  $B_i$  is the value of the  $i$ th bin. This metric converges to the standard definition of the relative  $L^2$  error as  $v \rightarrow 0$ .

These statistics cannot be used in general to find optimal values of  $\beta$  since they require knowledge of the analytic solution, and the algorithm must be run for a long time to build up a sufficiently large sample. However they can be used to assess the ability of other indicators to find the optimal proposal variances in a controlled setting. The following statistics can be used specifically for importance sampling algorithms.

**4.1.2. The variance of the weights.** Importance samplers assign a weight to each sample they produce based on a ratio of the posterior to the proposal at that point. Importance samplers are most efficient when the target is proportional to the proposal distribution. In this case the weights are all equal, and so the variance of the weights,  $\text{var}(w(y))$ , is zero. Hence, we would like to choose the value of  $\beta$  which minimises the variance of the weights,

$$\beta_{\text{var}}^* = \arg \min_{\beta} \text{var}(w(y)).$$

In our experience, the mean of the estimator  $\text{var}(w(y))$  is a smooth enough function of  $\beta$  that it can be used to tune the proposal variance during the burn-in phase of MCMC algorithms. However the variance of the estimator of the variance can be large, especially far away from the optimal value, so it can take a large number of iterations to calculate descent directions.

**4.1.3. The effective sample size.** The effective sample size,  $n_{\text{eff}}$ , can also be used to assess the efficiency of importance samplers. Ideally, in each iteration, we would like all  $M$  of our samples to provide us with new information about the posterior distribution. In practise, we cannot achieve a perfect effective sample size of  $M$ .

The effective sample size can be defined in the following way:

$$n_{\text{eff}} = \frac{\left( \sum_{i=1}^M w_i \right)^2}{\sum_{i=1}^M w_i^2} \approx \frac{M \mathbb{E}(w)^2}{\mathbb{E}(w^2)} = M \left( 1 - \frac{\text{var}(w)}{\mathbb{E}(w^2)} \right).$$

The second two expressions are true when  $M \rightarrow \infty$ . From the last expression we see that when the variance of the weights is zero,  $n_{\text{eff}} = M$ ; this is our ideal scenario. Maximising the effective sample size is equivalent to minimising the variance of the weights.

The statistic  $n_{\text{eff}}$  is easier to deal with than the variance of the weights, as it varies between 1 and  $M$  as opposed to the variance which can vary over many orders of magnitude. Therefore it is preferable as a means of tuning the scaling parameter.

In all of the numerics which follow, we use the effective sample size calculated at each iteration and averaged across the entire run to tune the scaling parameters. We do this because when we tune this parameter on the fly, we use the single-iteration average to estimate optimality. The optimal scaling parameter found using the global optimum of the variance of the weights is also included for comparison. Note that for the majority of iterations the value of  $n_{\text{eff}}$  is an overestimate of this statistic over a larger number of samples. Rare events which result in a large importance weight bring this statistic down, and care must be made not to overfit the proposal variance to  $n_{\text{eff}}$  on an iteration by iteration basis. This can be accounted for by increasing the variance slightly once the adaptive algorithm has arrived on a value using single iteration values of  $n_{\text{eff}}$ .

The effective sample size also has another useful property; if we imagine the algorithm in the burn-in phase, for example, we have  $M$  ensemble members in the tail of a Gaussian curve searching for the area of high density. If the ensemble members are evenly spaced, then the particle closest to the mean will have an exponentially higher weight assigned to it. The effective sample size ratio in this scenario will be close to 1. As the algorithm burns in, the ensemble populates the regions where the majority of the probability density lies, and the proposal distributions better represent the posterior distribution. This leads to smaller variance in the weights, and a bigger effective sample size. By this argument we can see that rising  $n_{\text{eff}}$  signals the end of the burn-in period.

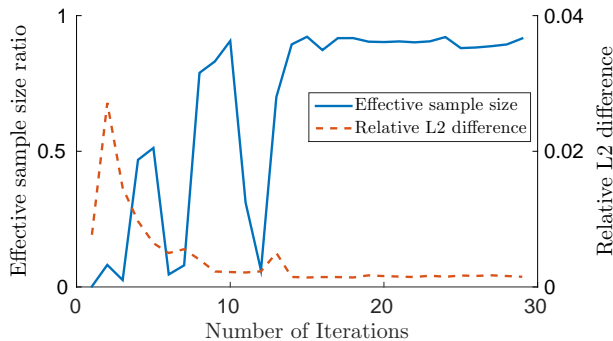
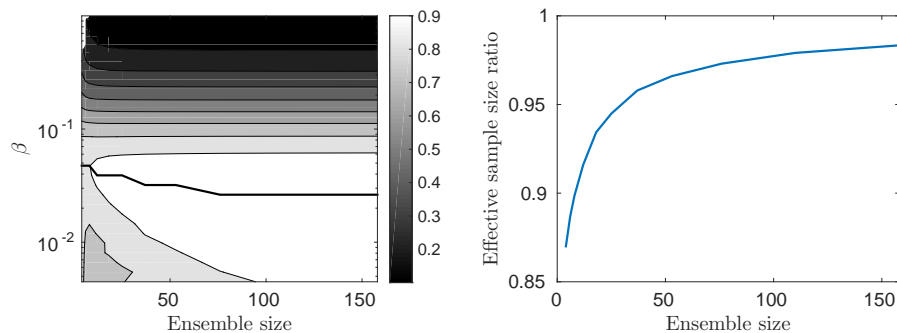


FIG. 4.1. The effective sample size ratio and relative  $L^2$  difference between the proposal and posterior distributions at each of the first 30 iterations. These numerics are taken from a simulation of the problem in Section 6.3 using the PAIS-Gamma algorithm.

Figure 4.1 shows that the effective sample size flattens out at the same time as the relative  $L^2$  difference between the posterior distribution and the proposal distribution stabilises close to its minimum. Detecting this stationarity allows us to automatically determine the end of the burn-in phase of the algorithm.

If we are to use the effective sample size ratio as an indicator of how to tune the proposal variance, we need to look at how it behaves in different situations. Here we look at how the statistic behaves as we vary the ensemble size,  $M$ . Figure 4.2 shows results for the Gaussian posterior discussed in Section 6.1 when using the PAIS algorithm with RW proposals, each with variance  $\beta^2 \in \mathbb{R}_{>0}$ . Figure 4.2 (a) shows that as the ensemble size increases, the scaling parameter which gives the optimal effective sample size ratio decreases. This is to be expected since if we are trying to approximate the posterior with a mixture distribution of a small number





(a) Contours showing optimal ranges of the scaling parameter. (b) The highest effective sample size achieved for each ensemble size.

FIG. 4.2. The behaviour of the effective sample size as the ensemble size increases. This analysis is from the Gaussian posterior in Section 6.1 using the PAIS-RW algorithm.

of Gaussians, the optimal variance will naturally be larger so that the whole of the significant regions are covered by the proposal distribution. As the number increases, the optimal variances decrease so that finer details in the posterior can be better represented in the proposal. Figure 4.2 (b) shows that as the ensemble size increases, the efficiency of the sampler also increases.

**4.2. Adaptive PAIS.** A popular approach for adaptive MCMC algorithms is to view the scaling parameter as a random variable which we can sample during the course of the MCMC iterations. However, it can be slow to converge to the optimal value, and we may need an uninformative prior for the scaling parameter. Alternatively, the parameter may be randomly sampled at various points during the evolution of the chain. This results in some iterations which make larger global moves in the state space between modes in the target distribution, and some which make local moves. Algorithms of this type do not converge to an optimal value of the scaling parameter. We choose to use a divide and conquer scheme which optimises the effective sample size (or any other diagnostic). Some more sophisticated examples are described in [23] and [13].

From here on in, we will use the random walk proposal,

$$y = x_{i-1} + \beta\omega_{i-1}, \quad \omega_{i-1} \sim \mathcal{N}(0, \Sigma) \quad \text{i.i.d.} \quad (4.3)$$

where  $x_{i-1}$  is our current state,  $y$  is our proposed state and  $\Sigma$  is a covariance operator, which could come from the prior distribution. In the numerics section we will refer to  $\beta$  as the scaling parameter.

Using an adaptive strategy, we calculate a sequence  $\{\beta^{(k)}\}_{k=1}$  which converges roughly to the optimal scaling parameter  $\beta^*$ , resulting in the optimal transition density  $\chi$  for our MCMC algorithm. This optimal value will differ depending on the criterion we are optimising. We must choose some sequence of iterations,  $\{n_k\}_{k=1}$ , at which to update  $\beta$ , and due to the constraints on adaptive MCMC algorithms [22, 23], these  $n_k$  must grow exponentially further apart. This same adaptive approach can also be applied to the trivially parallelised MCMC algorithms to adaptively calculate their optimal scaling parameter  $\beta^*$ .

**5. Approximate Multinomial Resampling.** Although the ETPF is optimal in terms of preserving statistics of the sample, it can also become quite costly as the number of ensemble members is increased. It is arguable that in the context of PAIS, we do not require this degree of accuracy, and that a faster more approximate method for resampling could be employed. One approach would be to use the bootstrap resampler, which simply takes the  $M$  ensemble members' weights and constructs a multinomial distribution, from which  $M$  samples are drawn. This is essentially the cheapest resampling algorithm that one could construct. However it too has some drawbacks. The algorithm is random, and as such it is possible for all of the ensemble members in a particular region not to be sampled. This could be particularly problematic when attempting to sample from a multimodal distribution, where it might take a long time to find one of the modes again. The bootstrap filter is also not guaranteed to preserve the mean of the weighted sample, unlike the ETPF. Ideally, we would like to use a resampling algorithm which is not prohibitively costly for moderately or large sized ensembles, which preserves the mean of the samples, and which makes it much harder for the new samples to forget a significant region in the density. This motivates the following algorithm, which we refer to as approximate multinomial resampling (AMR).

Instead of sampling  $M$  times from an  $M$ -dimensional multinomial distribution as is the case with the bootstrap algorithm, we sample once each from  $M$  different multinomials. Suppose that we have  $M$  samples  $y_n$  with weights  $w_n$ . The multinomial sampled from in the bootstrap filter has a vector of probabilities given by:

$$\frac{1}{\sum w_n} [w_1, w_2, \dots, w_M] = \bar{\mathbf{w}},$$

with associated states  $y_n$ . We wish to find  $M$  vectors  $\{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_M\} \subset \mathbb{R}_{\geq 0}^M$  such that  $\frac{1}{M} \sum \mathbf{p}_i = \bar{\mathbf{w}}$ . The AMR is then given by a sample from each of the multinomials defined by the vectors  $\mathbf{p}_i = [p_{i,1}, p_{i,2}, \dots, p_{i,M}]$  with associated states  $\mathbf{y}_i$ . Alternatively, as with the ETPF, a deterministic sample can be chosen by picking each sample to be equal to the mean value of each of these multinomial distributions, i.e. each new sample  $\hat{x}_i$  is given by:

$$\hat{x}_i = \sum p_{i,j} x_j, \quad i \in \{1, 2, \dots, M\}. \quad (5.1)$$

The resulting sample has several properties which are advantageous in the context of being used with the PAIS algorithm. Firstly, we have effectively chopped up the multinomial distribution used in the bootstrap filter into  $M$  pieces, and we can guarantee that exactly one sample will be taken from each section. This leads to a much smaller chance of losing entire modes in the density, if each of the sub-multinomials is picked in an appropriate fashion. Secondly, if we do not make a random sample for each multinomial with probability vector  $\mathbf{p}_i$  but instead take the mean of the multinomial to be the sample, this algorithm preserves the mean of the sample exactly. Lastly, as we will see shortly, this algorithm is significantly less computationally intensive than the ETPF.

There are of course infinitely many different ways that one could use to split the original multinomial up into  $M$  parts, some of which will be far from optimal. The method that we have chosen is loosely based on the idea of optimal transport. We search out states with the largest weights, and choose a cluster around these points based on the closest states geographically. This method is not optimal since once

most of the clusters have been selected the remaining states may be spread across the parameter space.

```

 $\mathbf{z} = M\bar{\mathbf{w}}$ 
for  $i = 1, 2, \dots, M$  do
     $J = \arg \max_j z_j$ 
     $p_{i,J} = \min\{1, z_J\}$ 
     $z_J = z_J - p_{i,J}$ 
    while  $\sum_j p_{i,j} < 1$  do
         $K = \arg \min_{k \in \{k | z_k > 0\}} \|y_J - y_k\|$ 
         $p_{i,K} = \min\{1 - \sum_j p_{i,j}, z_K\}$ 
         $z_K = z_K - p_{i,K}$ 
    end while
     $x_i = \sum_k p_{i,k} y_k$ 
end for

```

TABLE 5.1

*The approximate multinomial resampler (AMR) algorithm.*

Table 5.1 describes the basis of the algorithm with deterministic resampling, using the means of each of the sub-multinomials as the new samples. This resampler was designed with the aims of being numerically cheaper than the ETPF, and more accurate than straight multinomial resampling. Therefore we now present numerical examples which demonstrate this.

To test the accuracy and speed of the three resamplers (ETPF, bootstrap and AMR), we drew a sample of size  $M$  from the proposal distribution  $\mathcal{N}(1, 2)$ . Importance weights were assigned, based on a target distribution of  $\mathcal{N}(2, 3)$ . The statistics of the resampled outputs were compared with the original weighted samples. Figure 5.1 (a)-(c) show how the relative errors in the first three moments of the samples changes with ensemble size  $M$  for the three different samplers. As expected, the AMR lies somewhere between the high accuracy of the ETPF and the less accurate bootstrap resampling. Note that only the error for the bootstrap multinomial sampler is presented for the first moment since both the ETPF and the AMR preserve the mean of the original weighted samples up to machine precision. Figure 5.1 (d) shows how the computational cost, measured in seconds, scales with the ensemble size for the three different methods. These results demonstrate that the AMR behaves how we wish, and importantly ensures that exactly one sample of the output will lie in each region with weights up to  $\frac{1}{M}$  of the total.

We will use the AMR in the numerics in Section 6.3.3 where we have chosen to use a larger ensemble size. We do not claim that the AMR is the optimal choice within PAIS, but it does have favourable features, and demonstrates how different choices of resampler can affect the speed and accuracy of the PAIS algorithm.

## 6. Numerical Examples.

**6.1. Sampling from a one dimensional Gaussian distribution.** In this example we compare the naively parallelised RWMH algorithm with its PAIS variant, the PAIS-RW algorithm. The PAIS algorithm is implemented using the ETPF to perform the resampling step. We assess the performance of the PAIS algorithm using the relative  $L^2$  error defined in (4.2), as well as the relative error in the first moment (4.1).

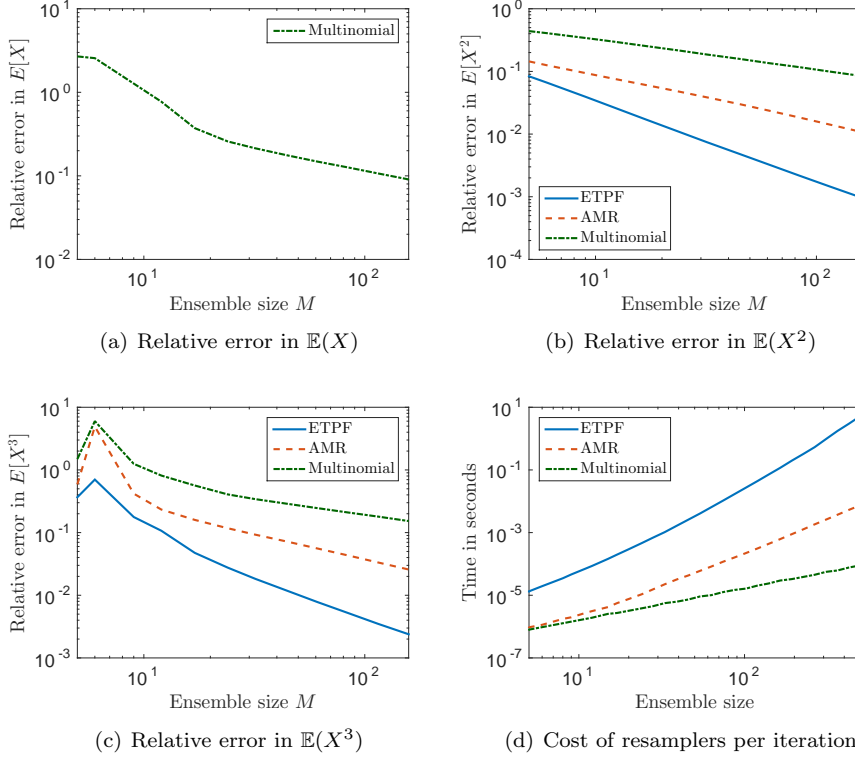


FIG. 5.1. Finding optimal values of  $\beta$  for the problem in Section 6.1. The setup is as in Section 6.1.2.

Since we are comparing against naively parallelised MH algorithms, we also need to decide which statistics  $T(\beta)$  provide the best criterions for obtaining the optimal scaling parameters. In the examples which follow, we have optimised the naively parallelised RWMH algorithm using the optimal acceptance rate  $\hat{\alpha} = 0.5$ . This value differs from the theoretical asymptotic value of 0.234 which applies in higher dimensions, but this higher acceptance rate is commonly used for one dimensional Gaussian posteriors [24]. To find the optimal scaling parameter we minimise the statistic

$$T_{\text{MH}}(\beta) = \left| \frac{N_{\text{acc}}(\beta)}{N_{\text{total}}} - \hat{\alpha} \right|,$$

where  $N_{\text{acc}}(\beta)$  is the number of accepted moves and  $N_{\text{total}}$  is the total number of samples produced. For the PAIS algorithm, we maximise the effective sample size as discussed in Section 4.1.3.

**6.1.1. Target distribution.** Consider the simple case of a linear observation operator  $\mathcal{G}(x) = x$ , where the prior on  $x$  and the observational noise follow Gaussian distributions. Then, following (2.2), the Gaussian posterior has the form

$$\text{law}(\mu_D) = \pi(x|D) \propto \exp\left(-\frac{1}{2}\|x - D\|_{\sigma^2}^2 - \frac{1}{2}\|x\|_{\tau^2}^2\right), \quad (6.1)$$

where  $\sigma^2$  and  $\tau^2$  are the variances of the observational noise and prior distributions respectively. In the numerics which follow, we choose  $\tau^2 = 0.01$  and  $\sigma^2 = 0.01$ , and we observe  $x_{\text{ref}} = 4$  noisily such that

$$D = \mathcal{G}(x) + \eta \sim \mathcal{N}(\mathcal{G}(x_{\text{ref}}), \sigma^2) = \mathcal{N}(4, 0.01).$$

These values result in a posterior density in which the vast majority of the density is out in the tails of the prior distribution. The Kullback-Leibler (KL) divergence, which gives us a measure of how different the prior and posterior are, is  $D_{KL}(\mu_D || \mu_0) = 4.67$  for this problem. A KL divergence of zero indicates that two distributions are identical almost everywhere.

**6.1.2. Numerical implementation.** In each of the following simulations, we perform three tasks. First we calculate the optimal value of  $\beta$  by optimising the statistics described in Section 4.1. We then run the algorithms with optimal parameters to calculate and compare the convergence rates. Finally, we implement the adaptive algorithms described in Section 4.2 and compare the convergence rates of these algorithms with the nonadaptive algorithms.

**(1) Finding the optimal parameters:** To find the optimal parameters we choose 32 values of  $\beta$  evenly spaced on a log scale in the interval  $[10^{-5}, 2]$ . We run the PAIS-RW and RWMH algorithms for one million iterations, each with an ensemble size of  $M = 50$ . We took 32 repeats of both algorithms and then used the geometric means of the sample statistics to find the optimal parameters.

**(2) Measuring convergence of nonadaptive algorithms:** We run the algorithms in Section 6.1 for one million iterations, again with  $M = 50$ . The simulations are repeated 32 times using the optimal parameters found in (1). The performance of the algorithms is judged by the convergence of the relative  $L^2$  error statistic in (4.2).

**(3) Measuring convergence of adaptive algorithms:** We run the adaptive algorithms under the same conditions as the nonadaptive algorithms, and again use the relative  $L^2$  error to compare efficiency. The adaptive algorithms are initialised with  $\beta^{(1)} = 1$ .

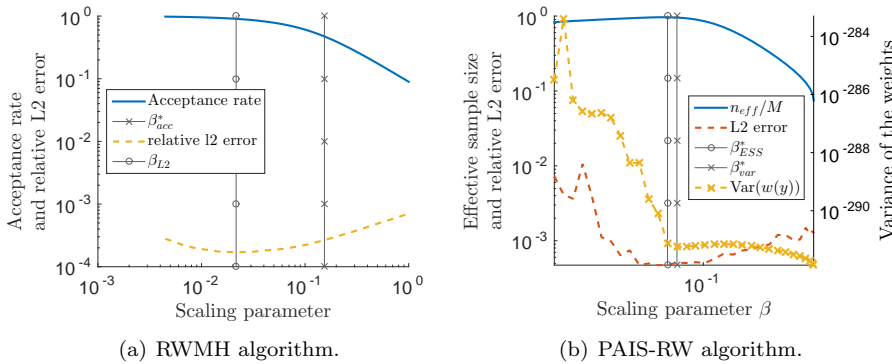


FIG. 6.1. Finding optimal values of  $\beta$  for the problem in Section 6.1. The setup is as in Section 6.1.2. Resampling is performed using the ETPF.

**6.1.3. Optimal values of  $\beta$ .** Figure 6.1 (a) shows the two values of  $\beta$  which are optimal according to the acceptance rate and relative  $L^2$  error criteria for the RWMH algorithm. The smaller estimate comes from the relative  $L^2$  error, and the larger from

Statistic	RWMH
$\beta_{L^2}^*$	2.1e-2
$\beta_{\%}^*$	1.5e-1
Acceptance Rate ( $\beta_{L^2}^*$ )	9.0e-1
Acceptance Rate ( $\beta_{\%}^*$ )	5.0e-1

Statistic	PAIS-RW
$\beta_{\text{eff}}^*$	4.7e-2
$\beta_{\text{var}(w(y))}^*$	5.8e-2
$\beta_{L^2}^*$	3.9e-2

TABLE 6.1

Optimal values of  $\beta$  summarised from Figure 6.1. Statistics calculated as described in Section 4.1. The values  $\beta_{L^2}^*$  and  $\beta_{\%}^*$  are the optimal scaling parameters found by optimising the relative  $L^2$  errors and acceptance rate respectively. Similarly  $\beta_{\text{eff}}^*$  and  $\beta_{\text{var}(w(y))}^*$  optimise the effective sample size and variance of the weights statistics.

the acceptance rate. The results in Figure 6.1 are summarised in Table 6.1. Since in general we cannot calculate the relative  $L^2$  error, we must optimise the algorithm using the acceptance rate. From the relative  $L^2$  error curve we can see that the minimum is very wide and despite the optimal values being very different there is not a large difference in the convergence rate.

Figure 6.1 (b) shows the effective sample size ratio compared to the error analysis and the variance of the weights. The relative  $L^2$  error graph is noisy, but it is clear that the maximum in the effective sample size and the minimum in the variance of the weights are both close to the minimum in the relative  $L^2$  error. Due to this we say that the estimate of the effective sample size found by averaging the statistic over each iteration is a good indicator for the optimal scaling parameter. In general this indicator overestimates the value of  $n_{\text{eff}}$  found by using the entire sample.

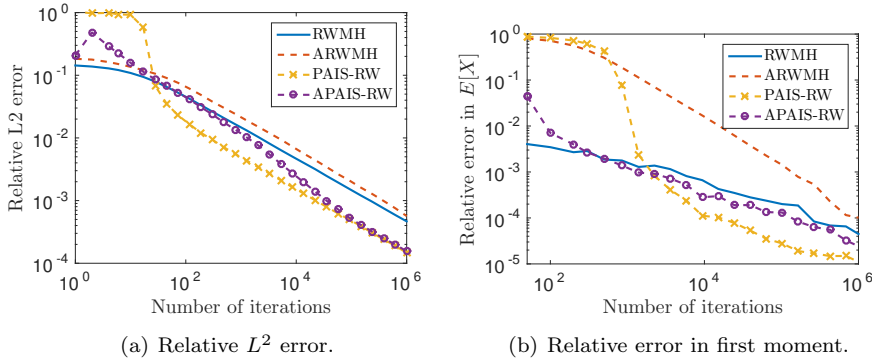


FIG. 6.2. Error analysis for the (A)RWMH and (A)PAIS-RW algorithms. The setup is as in Section 6.1.2 (2, 3). Resampling is performed using the ETPF.

**6.1.4. Convergence of RWMH vs PAIS-RW.** Figure 6.2 shows that the PAIS-RW algorithm converges to the posterior distribution significantly faster than the RWMH algorithm, in both  $L^2$  error and relative error in the moments. A description of the speed up attained by this algorithm is given in Section 6.2.4.

Both adaptive algorithms are run with initial values of  $\beta = 1$ . Figure 6.2 shows that after an initial burn-in period the APAIS-RW algorithm catches up to the PAIS-RW algorithm, and by the end of the simulation window is matching its performance. The ARWMH algorithm does not perform quite as well, this is possibly due to the fact

that the acceptance rate cost function is not particularly smooth at the optimal value making it difficult to minimise.

**6.1.5. Scaling of the PAIS algorithm with ensemble size.** Throughout this example, we use an ensemble size  $M = 50$ , but it is interesting to see how the PAIS algorithm scales when we increase the ensemble size, and if there is some limit below which the algorithm fails. We implement the problem in Section 6.1, using the RWMH and PAIS-RW algorithms with ensemble sizes in the interval  $M \in [1, 160]$ .

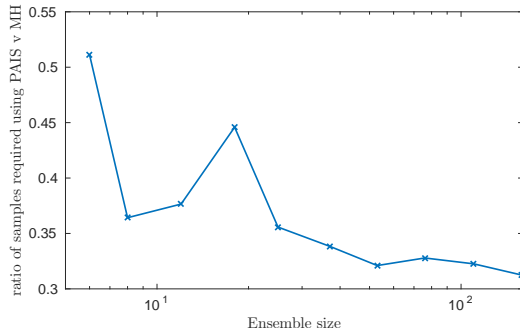


FIG. 6.3. *Ratio of PAIS-RW samples required to reach the same tolerance as the RWMH algorithm.*

Figure 6.3 was produced using the method of finding optimal  $\beta$  described in Section 6.1.2 (1), then running 32 repeats at each ensemble size. The convergence rates are then found by regressing through the data. The graph is still very noisy but demonstrates that increasing the ensemble size continues to reduce the number of iterations required in comparison with naively parallelised MH. The decreasing trend indicates superlinear improvement of PAIS with respect to ensemble size, in terms of the number of iterations required, which is a demonstration of our belief that parallelism of MCMC should give us added value over and above that provided by naive parallelism. This decrease is linked to the increasing effective sample size shown in Figure 4.2 (b).

**6.2. Sampling from Bimodal Distributions.** In this section we investigate the behaviour of the PAIS algorithm when applied to bimodal problems. MH methods can struggle with multimodal problems, particularly where switches between the modes are rare, resulting in incorrectly proportioned modes in the histograms. This example demonstrates that the PAIS algorithm redistributes chains to new modes as they are found. This means that we expect the number of chains in a mode to be approximately proportional to the probability density in that mode. As a result, reconstructed posteriors with disproportional modes, as is familiar with the MH algorithms, are not produced.

**6.2.1. Target Distribution.** We look at an ‘easy’ problem,  $B_1$ , which has a KL divergence of 0.880, and a ‘harder’ problem,  $B_2$ , which has a KL divergence of 3.647. Problem  $B_1$  has two modes which are not too far apart. In  $B_2$  we increase the distance between the two modes which has the effect of increasing the expected number of iterations that it takes for a MCMC chain to jump between modes. These posteriors are shown in Figure 6.4.

The following setup is the same for both problems. We consider a non-linear observation operator  $\mathcal{G}(x) = x^2$ , and assign the prior  $x \sim \mu_0 = \mathcal{N}(0, \tau^2 = 0.25)$ . We assume that a noisy reading,  $D$ , is taken according to  $D = \mathcal{G}(x_{\text{ref}}) + \varepsilon$ , where  $\varepsilon \sim \mu_\varepsilon = \mathcal{N}(0, \sigma^2 = 0.1)$ . This results in the non-Gaussian posterior

$$\pi(x|D) \propto \exp\left(-\frac{1}{2\sigma^2}\|x^2 - D\|^2 - \frac{1}{2\tau^2}\|x\|^2\right).$$

To create the ‘easy’ problem we say that the true value of  $\mathcal{G}(x_{\text{ref}}) = 0.75$ , and the ‘hard’ problem is generated using  $\mathcal{G}(x_{\text{ref}}) = 2$ . In the numerics which follow we draw noise from  $\mu_\varepsilon$  to generate our data point.

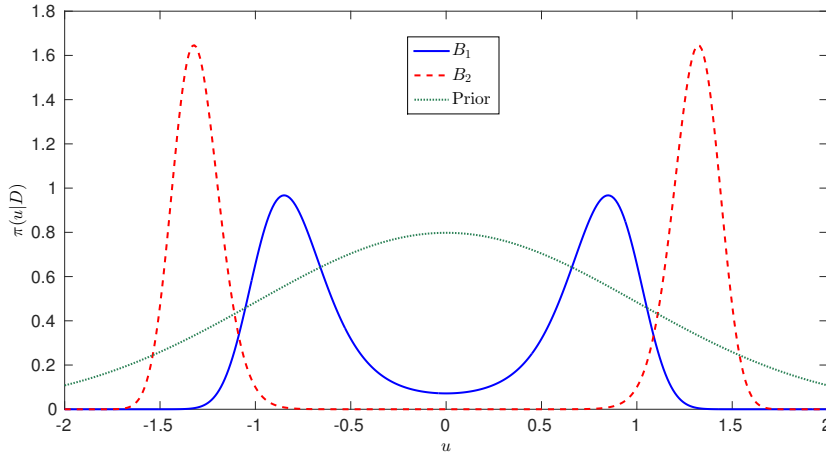


FIG. 6.4. Posterior distributions for problems  $B_1$  and  $B_2$  as described in Section 6.2.1.

**6.2.2. Calculating values of Optimal  $\beta^*$ .** Calculating the optimal values of the scaling parameters for this problem is similar to the previous example; we check only the acceptance rate to find the optimal values for RWMH and we use the effective sample size to find the optimal values for PAIS-RW. Table 6.2 gives the optimal values of  $\beta$  for both problems. The subscript on  $\beta$  refers to the criterion which has been optimised.

Algorithm	$\beta_{\text{acc}}^*$	$\beta_{\text{eff}}^*$
RWMH	4.8e-1	-
PAIS-RW	-	1.0e-1

Algorithm	$\beta_{\text{acc}}^*$	$\beta_{\text{eff}}^*$	$\beta_{\text{L2}}^*$
RWMH	2.3e-1	-	9.3e-1
PAIS-RW	-	5.1e-2	1.3e-1

TABLE 6.2  
Optimal values of  $\beta$  for  $B_1$  (left) and  $B_2$  (right).

It is relatively simple to find optimal scaling parameters for problem  $B_1$ . These values are given in Table 6.2 (left). However problem  $B_2$  is much harder as transitions between the modes are extremely unlikely for the standard RWMH algorithm. This means that we need to consider the convergence on two levels; we should consider the algorithm’s ability to find both the modes, and also whether it can sample them in the correct proportions.



To get correctly proportioned modes with the RWMH algorithm it is important that the chains can transition between the modes frequently, which means that  $\beta$  must be large. However, this leads to a lower acceptance rate, and so we sacrifice convergence locally. For this reason, the RWMH algorithm is very slow to converge for problems of this type.

We can achieve these two regimes in RWMH by tuning  $\beta$  using the acceptance rate for local convergence, and by  $L^2$  error for global convergence. Similarly in PAIS-RW we can use the effective sample size for local convergence, and the  $L^2$  error for global convergence.

From Table 6.2 (right) we see that there is a large difference between the optimal value of  $\beta$  for each regime in RWMH, and so will result in inefficient sampling. The PAIS-RW algorithm manages to sample the local detail and the large scale behaviour with similar values of  $\beta^*$ ; a clear advantage to using this algorithm for this problem.

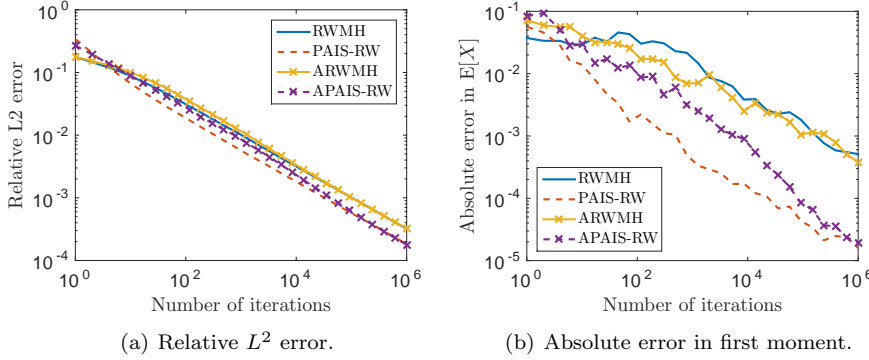


FIG. 6.5. Convergence of the PAIS-RW and RWMH algorithms for Problem  $B_1$ . Set up described in Section 6.1.2. Resampling is performed using the ETPF.

**6.2.3. Convergence of RWMH vs PAIS-RW.** As in the Gaussian example we see a significant speed up with the PAIS-RW algorithm for problem  $B_1$ . Figure 6.5 shows the adaptive and nonadaptive convergence rates. We can see that the adaptive algorithms compare closely with the respective nonadaptive algorithms and the improvement PAIS offers remains significant.

For  $B_2$  the algorithms are run with the global optimal value of  $\beta^*$ , and with the local optimal value of  $\beta^*$ . Figure 6.6 (a) shows that the algorithms using the globally optimal  $\beta^*$  converge at the desired rate, whereas the algorithms optimised using the acceptance rate and effective sample size initially converge faster but at some point forget the location of one of the modes, causing the convergence to flatten out. Using the PAIS algorithm we can minimise the impact of forgotten modes by constantly allowing the algorithm to search them out. Since we have parallel chains, we can run PAIS-RW with the majority of chains using the locally optimal value of  $\beta$ , and one or two chains with a larger scaling parameter. These chains with larger scaling parameters act both as ‘scouts’ for new modes, and act to aid in the overdispersal of the proposal distribution. Figure 6.6 (b) shows the results of using 49 chains with the local optimal scaling parameter, and one chain with ten times the local optimal scaling parameter. We see that modes are not forgotten and the algorithm converges with the improvement we see from PAIS-RW in the other problems. Other methods of mode searches are described in [15].

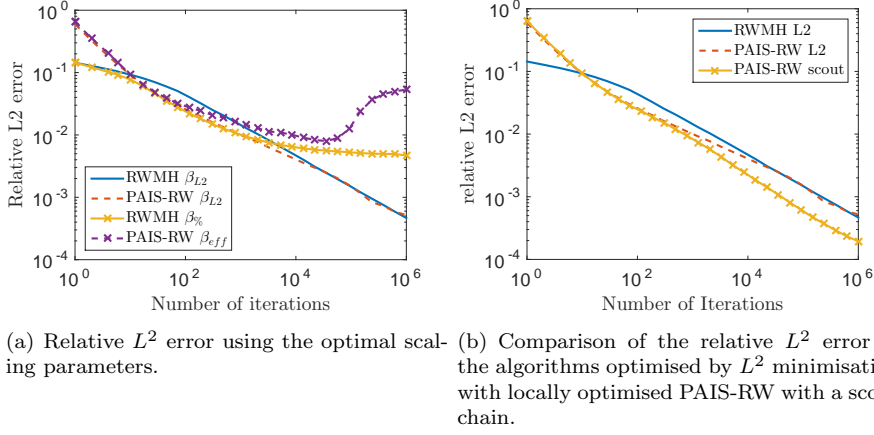


FIG. 6.6. Convergence of the PAIS-RW and RWMH algorithms for Problem  $B_2$ . Set up described in Section 6.1.2. Resampling is performed using the ETPF.

The adaptive algorithm can just as easily be applied to the PAIS algorithm with the ‘scout’ chains described in the previous paragraph. Since we need two equally sized sub-ensembles, we will use two groups of 24 ensembles with the same proposal distributions, and each group will also have a ‘scout’ chain with a scaling parameter ten times that of the rest of the group.

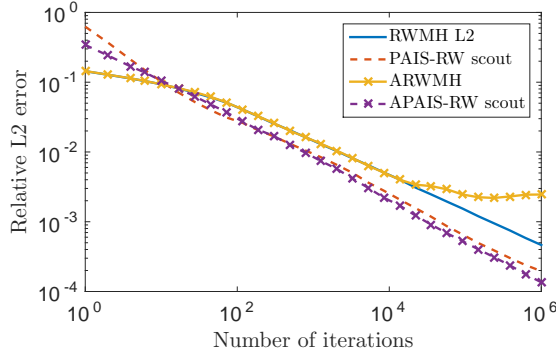


FIG. 6.7. Convergence of the relative  $L^2$  error for problem  $B_2$ , comparing the globally optimised nonadaptive algorithms with the locally optimised adaptive algorithms. The setup is as described in Section 6.1.2 (3). Resampling is performed using the ETPF.

Comparing the convergence of the adaptive algorithms against the nonadaptive algorithms in Figure 6.7 shows that the algorithms behave as expected. The adaptive RWMH algorithm tuned using the acceptance rate converges at the same rate as the  $L^2$  optimised algorithm until the scaling parameter gets small, and therefore switches between the modes are rare, and the relative heights of the modes are decided by the arbitrary proportion of chains which are in each mode at this point. The adaptive PAIS-RW with scout chains tuned to the effective sample size converges at about the same rate as the locally optimised nonadaptive algorithm also with scouts.

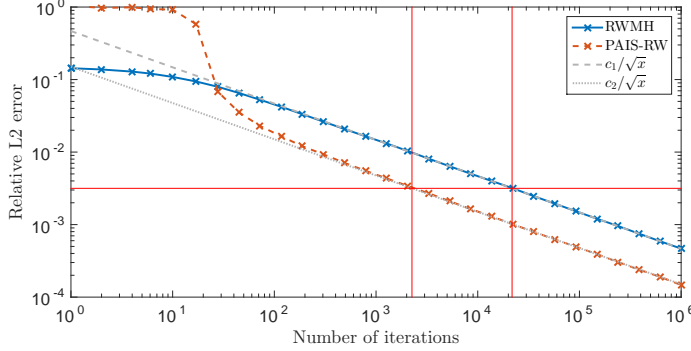


FIG. 6.8. Illustration of calculating the number of PAIS-pCNL iterations required to reach a tolerance of  $10^{-2.8}$  as a percentage of MH iterations. The relative  $L^2$  error graphs are from the Gaussian problem in Section 6.1.

**6.2.4. Calculating the Speed Up in Convergence.** The graphs in the previous section clearly show that the PAIS-RW algorithm converges faster than the RWMH algorithm when both are parallelised with the same ensemble size. We can calculate the number of iterations required to achieve a particular tolerance level in our solution for each algorithm and compare these to calculate a percentage saving. In Figure 6.8 we demonstrate our calculation of the savings. The constants  $c_1$  and  $c_2$  are found by regressing through the data with a fixed exponent of  $-1/2$  excluding the initial data points where the graph has not finished burning in.

A summary of the percentage of iterations required using the PAIS algorithm compared with the respective Metropolis-Hastings algorithms is given in Table 6.3. The blank entries correspond to occasions when either the MH algorithm or PAIS algorithm hasn't converged to the posterior distribution.

	Gaussian	$B_1$	$B_2$
RWMH	10%	32%	12.6% (scout)
pCN	-	32%	-
MALA	42%	36%	40%
pCNL	66%	56%	-

TABLE 6.3

Iterations for the PAIS algorithms required to achieve a desired tolerance as a percentage of the number of iterations required by the respective MH algorithms. The pCN and pCNL proposal distributions are taken from [6].

**6.2.5. A Useful Property of the PAIS Algorithm for Multimodal Distributions.** The biggest issue for the Metropolis-Hastings algorithms when sampling from a posterior such as the one in  $B_2$  is that it is unlikely that the correct ratio of chains will be maintained in each of the modes, and since there is no interaction between the chains, there is no way to remedy this problem. The PAIS algorithm tackles this problem with its resampling step. The algorithm uses its dynamic kernel to build up an approximation of the posterior at each iteration, and then compares this to the posterior distribution via the weights function. Any large discrepancy in the approximation will result in a large or small weight being assigned to the relevant chain, meaning the chain will either pull other chains towards it or be sucked towards

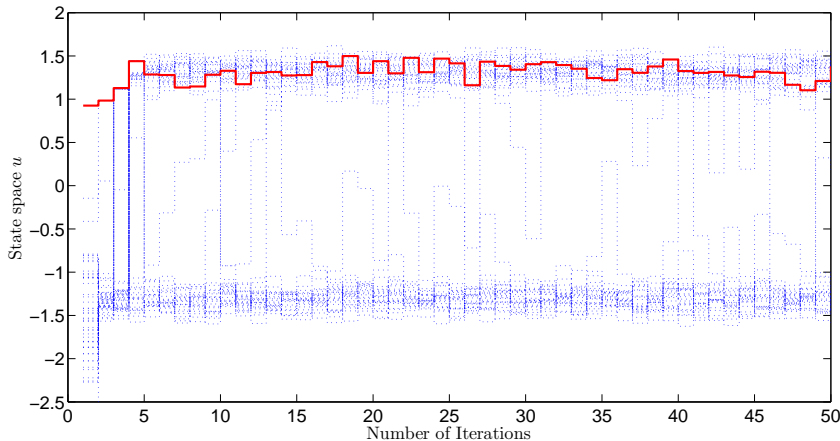


FIG. 6.9. This figure demonstrates the redistribution property of the PAIS algorithm. Initially there is one chain in the positive mode, and 49 chains in the negative mode.

a chain with a larger weight. In this way, the algorithm allows chains to ‘teleport’ to regions of the posterior which are in need of more exploration. Figure 6.9 shows Problem  $B_2$  with initially 1 chain in the positive mode, and 49 chains in the negative mode. It takes only a handful of iterations for the algorithm to balance out the chains into 25 chains in each mode. The chains switch modes without having to climb the energy gradient in the middle.

**6.3. Sampling from non-Gaussian bivariate distributions.** In this section we apply the PAIS algorithm to a more complicated posterior distribution. The field of biochemical kinetics gives rise to multiscale stochastic problems which remain a challenge both theoretically and computationally. Biochemical reactions occur in single cells between a number of chemical populations and the rates of these reactions can vary on vastly different timescales. It is these reaction rates which we are interested in finding descriptions for.

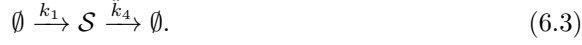
It is often possible to isolate which reactions are occurring more frequently (the fast reactions) and which are occurring less frequently (the slow reactions). The quasi-steady-state assumption (QSSA) is the assumption that the fast reactions converge in distribution on a timescale which is negligible with respect to the rate of occurrence of the slow reactions. This assumption allows us to approximate the dynamics of the slowly changing quantities in the system by assuming that the fast quantities are in equilibrium with respect to the fast reactions in isolation. This kind of model reduction can be used to approximate the likelihood in an inverse problem where we wish to recover the reaction parameters in the system.

Let us consider a simple example by introducing the following simple chemical system involving two chemical species  $S_1$  and  $S_2$ :



Each arrow represents a reaction from a reactant to a product, with some rate constant  $k_i$ , and where the rates of the reactions are assumed to follow mass action kinetics. We denote the concentration of species  $S_i$  by  $X_i$ . We assume that we are

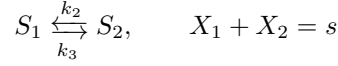
in a parameter regime such that the reactions  $S_1 \rightarrow S_2$  and  $S_2 \rightarrow S_1$  occur much more frequently than the other reactions. Notice that both chemical species are involved in fast reactions. However, the quantity  $\mathcal{S} = X_1 + X_2$  is conserved by both of the fast reactions, and as such, this is the slowly changing quantity in this system. The effective dynamics of  $\mathcal{S}$  can be represented as follows.



Here, the new reaction rate  $\hat{k}_4$  is approximated through application of the QSSA to be

$$\hat{k}_4(s) = \mathbb{E}[k_4 X_2 | \mathcal{S} = s] = \frac{k_2 k_4 \mathcal{S}}{k_2 + k_3}.$$

The value of  $\mathbb{E}[k_4 X_2 | \mathcal{S} = s]$  is approximated by finding the steady state of the ODE representing the fast subsystem of reactions:



If we assume that we know the rate constants  $k_1$  and  $k_4$ , and observe the system in (6.3), we indirectly observe the rates  $k_2$  and  $k_3$  through the effective rate  $\hat{k}_4$  of the degradation of  $\mathcal{S}$ . Our observations are uninformative about these reaction rates, as there are surfaces in parameter space along which the effective rate  $\hat{k}_4$  is invariant, leading to a highly ill-posed inverse problem. Making the assumption that the errors in our observations of the value of  $\mathcal{S}$  are Gamma distributed with some variance  $\sigma^2$ , this results in a long and thin posterior distribution on  $k_2$  and  $k_3$ . This type of problem is notoriously difficult to sample from using standard MH algorithms, as the algorithms quickly find a point on this manifold on which  $\hat{k}_4$  is invariant, but exploration along its length is slow.

**6.3.1. Target Distribution.** We now formalise the posterior of interest. We look for a distribution over the parameter  $\mathbf{k} = (k_2, k_3)^T$ , given that we know  $k_1 = 100$  and  $k_4 = 1$ . We generate our data by making ten observations of the system at  $t_i = 2, 4, \dots, 20$ , here simulated by solving the full system, (6.2), governed by the differential equations

$$\begin{aligned} \frac{dS_1}{dt} &= k_1 - k_2 X_1(t) + k_3 X_2(t), \\ \frac{dS_2}{dt} &= k_2 X_1(t) - (k_3 + k_4) X_2(t), \end{aligned}$$

with initial conditions  $X_1(0) = X_2(0) = 0$ , and parameter values  $\mathbf{k} = (50, 100)^T$ . We then add noise taken from a Gamma distribution with variance  $\sigma^2 = 225$  and centred at each  $\mathcal{S}(t_i) = X_1(t_i) + X_2(t_i)$ .

In our modelling we use the QSSA to simplify this system of differential equations into the one dimensional system based on (6.3),

$$\frac{d\mathcal{S}}{dt} = k_1 - \frac{k_2 k_4}{k_2 + k_3} \mathcal{S}(t).$$

The observation operator,  $\mathcal{G} : (\mathbf{k}, t) \mapsto \mathcal{S}(t)$ , maps from parameter space onto population space at a time  $t$ . This means that for the  $i$ th observation we assume,

$$D_i \sim \text{Gamma}(\alpha_i, \beta_i),$$

where

$$\alpha_i = \frac{\mathcal{G}(\mathbf{k}, t_i)^2}{\sigma^2}, \quad \beta_i = \frac{\mathcal{G}(\mathbf{k}, t_i)}{\sigma^2}, \quad i = 1, \dots, 10.$$

These are parameters required to give us a distribution with mean  $\mathcal{G}(\mathbf{k})$  and variance  $\sigma^2$ . We assign Gamma priors to  $\mathbf{k}$  with mean  $\alpha_0/\beta_0 = 75$  and variance  $\alpha_0/\beta_0^2 = 100$  in both coordinates, resulting in the posterior

$$\pi(\mathbf{k}|\mathbf{D}) \propto \left[ \prod_{i=1}^{10} \text{Gamma}(D_i; \alpha_i, \beta_i) \right] \text{Gamma}(k_2; \alpha_0, \beta_0) \text{Gamma}(k_3; \alpha_0, \beta_0).$$

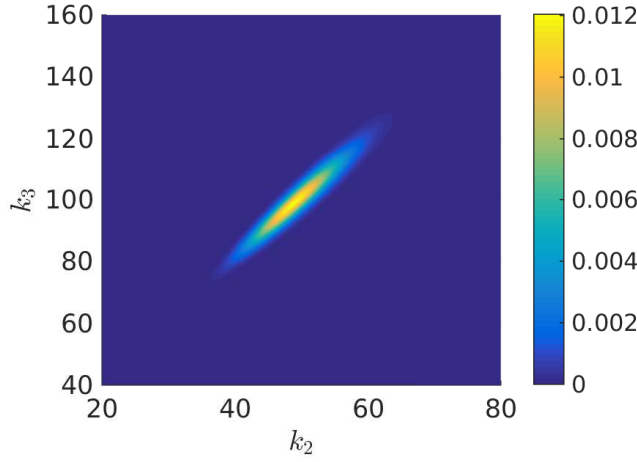


FIG. 6.10. The posterior distribution on the parameters  $k_2$  and  $k_3$  given the data and priors described in Section 6.3.1.

Figure 6.10 presents a visualisation of the posterior distribution for this problem, found by exhaustive MH simulation.

**6.3.2. Implementation.** For this problem there is no analytic form for the normalisation constant, and numerical methods implemented in MATLAB and Mathematica have proven to be unreliable. To demonstrate the convergence of the algorithm, we first run a MH algorithm for much longer than we normally would ( $8 \times 10^{10}$  samples), and consider the histogram produced to be a well-converged approximation of the true posterior distribution. We then perform our usual simulations with the PAIS and MH algorithms to compare the rate at which the histograms converge to the posterior over the same mesh.

In this problem we modify our PAIS algorithm to use a mixture of Gamma distributions in the proposal distribution instead of a Gaussian mixture. The PAIS-Gamma and MH-Gamma algorithms use a Gamma proposal distribution with mean centred at the previous state,

$$y \sim \text{Gamma}(\cdot; \alpha^*, \beta^*) \quad \text{where} \quad \frac{\alpha^*}{\beta^*} = x \quad \text{and} \quad \frac{\alpha^*}{(\beta^*)^2} = \beta^2.$$

Similarly the PAIS-LA algorithm uses the Langevin proposal distribution from the

standard MALA algorithm with a perturbation taken from a Gamma distribution,

$$y \sim \text{Gamma}(\cdot; \alpha^*, \beta^*), \quad \text{where} \quad \frac{\alpha^*}{\beta^*} = x + \frac{1}{2}\beta^2 \nabla \log \pi(X) \quad \text{and} \quad \frac{\alpha^*}{(\beta^*)^2} = \beta^2.$$

This ensures that we only propose positive parameters. Having a heavier right tail in the proposal distribution also means that the posterior is absolutely continuous with respect to the proposal. When this is not the case, the weight function can tend to zero or infinity when  $k_i \rightarrow \infty$ , which results in poor, very spiky approximations of the tails of the posterior distribution leading to slow convergence.

We now significantly increase the ensemble size we are using from  $M = 50$  to  $M = 2500$ . This allows us to build a better approximation of the posterior distribution for our proposals. Following the discussion in Section 5, it is clear that to keep the runtime of the resampler negligible compared with the calculation of the posterior we must switch from the ETPF resampler to the AMR algorithm.

**6.3.3. Convergence of PAIS-Gamma and PAIS-LA vs MH with Gamma proposals.** In the numerics which follow we have used the AMR algorithm to resample with an ensemble size of  $M = 2500$ . The method otherwise remains the same as in previous sections. We perform test runs to find the optimal scaling parameters for both Gamma proposals and MALA-type proposals. We then calculate the convergence rates of the algorithms by producing 50 million samples from the posterior with each algorithm, and repeat the simulation 32 times.

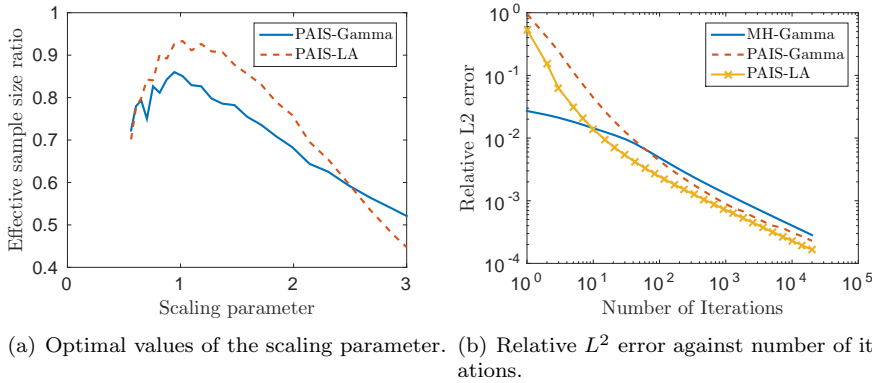


FIG. 6.11. Convergence of the PAIS-Gamma and PAIS-LA algorithms for the chemical system problem described in Section 6.3. Implementation described in Sections 6.3.2 and 6.3.3. Resampling is performed using the AMR scheme.

Algorithm	MH-Gamma	PAIS-Gamma	PAIS-LA
$\beta^*$	2.7e-0	9.4e-1	1.0e-0

TABLE 6.4

Optimal values of the scaling parameter. The MH algorithm is optimised using the acceptance rate, and the PAIS algorithms are optimised using the effective sample size.

In Figure 6.11 (a) we see that the two PAIS algorithms have similar optimal values of the scaling parameter, also shown in Table 6.4. However, the MALA-type proposal

achieves a higher effective sample size with the same ensemble size. Convergence after the first 50 million samples is shown in Figure 6.11 (b) and demonstrates that the PAIS algorithm converges faster than the MH algorithm. The slightly unstable convergence in the PAIS-Gamma algorithm is the effect of spikes appearing in the tails of the posterior distribution. This is also what has caused the slightly lower effective sample size.

**7. Discussion and Conclusions.** We have explored the application of parallelised MCMC algorithms in low dimensional inverse problems. We have demonstrated numerically that these algorithms converge faster than the analogous naively parallelised Metropolis-Hastings algorithms. Further experimentation with the Metropolis Adjusted Langevin Algorithm (MALA), preconditioned Crank-Nicolson (pCN), preconditioned Crank-Nicolson Langevin (pCNL) and Hamiltonian Monte Carlo (HMC) proposals has yielded similar results [25].

Importantly, we have compared the efficiency of our parallel scheme with a naive parallelisation of serial methods. Thus our increase in efficiency is over and above an  $M$ -fold increase, where  $M$  is the number of ensemble members used. Our approach demonstrates a better-than-linear speed-up with the number of ensemble members used.

The PAIS has a number of favourable features, for example the algorithm's ability to redistribute, through the resampling regime, the ensemble members to regions which require more exploration. This allows the method to be used to sample from complex multimodal distribution.

Another strength of the PAIS is that it can also be used with any MCMC proposal. There are a growing number of increasing sophisticated MCMC algorithms (HMC, Riemann manifold MCMC etc) which could be incorporated into this framework, leading to even more efficient algorithms, and this is another opportunity for future work.

One limitation of the PAIS approach as described above is that a direct solver of the ETPF problem (such as FastEMD [19]) has computational cost  $\mathcal{O}(M^3 \log M)$ , where  $M$  is the number of particles in the ensemble. As such, we introduced a more approximate resampler the approximate multinomial resampler, which allows us to push the approach to the limit with much larger ensemble sizes. The PAIS framework is very flexible in terms of being able to use any combination of proposal distributions and resampling algorithms that one wishes.

## REFERENCES

- [1] A. BESKOS, F. PINSKI, J. SANZ-SERNA, AND A. STUART, *Hybrid monte carlo on hilbert spaces*, Stochastic Processes and their Applications, 121 (2011), pp. 2201–2230.
- [2] T. BUI-THANH AND M. GIROLAMI, *Solving large-scale pde-constrained bayesian inverse problems with riemann manifold hamiltonian monte carlo*, Inverse Problems, 30 (2014), p. 114014.
- [3] B. CALDERHEAD, *A general construction for parallelizing metropolis- hastings algorithms*, Proceedings of the National Academy of Sciences, 111 (2014), pp. 17408–17413.
- [4] C. COTTER AND S. REICH, *Ensemble filter techniques for intermittent data assimilation-a survey*, in Large Scale Inverse Problems. Computational Methods and Applications in the Earth Sciences, Walter de Gruyter, Berlin, 2012, pp. 91–134.
- [5] S. COTTER, M. DASHTI, J. ROBINSON, AND A. STUART, *Bayesian inverse problems for functions and applications to fluid mechanics*, Inverse Problems, 25 (2009), p. 115008.
- [6] S. COTTER, G. ROBERTS, A. STUART, AND D. WHITE, *MCMC methods for functions: modifying old algorithms to make them faster*, Statistical Science, 28 (2013), pp. 424–446.
- [7] T. EL MOSELHY AND Y. MARZOUK, *Bayesian inference with optimal maps*, Journal of Computational Physics, 231 (2012), pp. 7815–7850.



- [8] G. EVENSEN, *Sequential data assimilation with a nonlinear quasi-geostrophic model using Monte Carlo methods to forecast error statistics*, Journal of Geophysical Research: Oceans (1978–2012), 99 (1994), pp. 10143–10162.
- [9] A. GELMAN AND D. B. RUBIN, *Inference from Iterative Simulation Using Multiple Sequences*, Statistical Science, 7 (1992), pp. 457–472.
- [10] M. GIROLAMI AND B. CALDERHEAD, *Riemann manifold langevin and hamiltonian monte carlo methods*, Journal of the Royal Statistical Society: Series B (Statistical Methodology), 73 (2011), pp. 123–214.
- [11] N. GORDON, D. SALMOND, AND A. SMITH, *Novel approach to nonlinear/non-Gaussian Bayesian state estimation*, in IEE Proceedings F (Radar and Signal Processing), vol. 140, IET, 1993, pp. 107–113.
- [12] W. HASTINGS, *Monte Carlo sampling methods using Markov chains and their applications*, Biometrika, 57 (1970), pp. 97–109.
- [13] C. JI AND S. C. SCHMIDLER, *Adaptive Markov Chain Monte Carlo for Bayesian Variable Selection*, Journal of Computational and Graphical Statistics, 22 (2013), pp. 708–728.
- [14] R. KALMAN, *A new approach to linear filtering and prediction problems*, Journal of Fluids Engineering, 82 (1960), pp. 35–45.
- [15] S. LAN, J. STREETS, AND B. SHAHBABA, *Wormhole Hamiltonian Monte Carlo*, ArXiv e-prints, (2013).
- [16] J. LIU, F. LIANG, AND W. WONG, *The multiple-try method and local optimization in metropolis sampling*, Journal of the American Statistical Association, 95 (2000), pp. 121–134.
- [17] G. MOORE ET AL., *Cramming more components onto integrated circuits*, Proceedings of the IEEE, 86 (1998), pp. 82–85.
- [18] R. NEAL, *Mcmc using ensembles of states for problems with fast and slow variables such as gaussian process regression*, arXiv preprint arXiv:1101.0387, (2011).
- [19] O. PELE AND M. WERMAN, *Fast and robust Earth mover’s distances*, in Computer Vision, 2009 IEEE 12th International Conference on, IEEE, 2009, pp. 460–467.
- [20] S. REICH, *A nonparametric ensemble transform method for Bayesian inference*, SIAM Journal on Scientific Computing, 35 (2013), pp. A2013–A2024.
- [21] G. ROBERTS AND J. ROSENTHAL, *Optimal scaling for various Metropolis-Hastings algorithms*, Statistical science, 16 (2001), pp. 351–367.
- [22] ———, *Coupling and ergodicity of adaptive Markov chain Monte Carlo algorithms*, Journal of Applied Probability, 44 (2007), pp. 458–475.
- [23] ———, *Examples of adaptive MCMC*, Journal of Computational and Graphical Statistics, 18 (2009), pp. 349–367.
- [24] J. ROSENTHAL, *Optimal proposal distributions and adaptive MCMC*, Handbook of Markov Chain Monte Carlo, (2011), pp. 93–112.
- [25] P. RUSSELL, *PhD Thesis*, PhD thesis, School of Mathematics, University of Manchester, 2017.
- [26] J. SEXTON AND D. WEINGARTEN, *Hamiltonian evolution for the hybrid monte carlo algorithm*, Nuclear Physics B, 380 (1992), pp. 665–677.
- [27] A. STUART, *Inverse problems: a Bayesian perspective*, Acta Numerica, 19 (2010), pp. 451–559.
- [28] C. VILLANI, *Topics in Optimal Transportation*, Graduate studies in mathematics, American Mathematical Society, 2003.
- [29] ———, *Optimal Transport: Old and New*, Grundlehren der mathematischen Wissenschaften, Springer Berlin Heidelberg, 2008.